

1. (10 points) What is the distance between two parallel hyperplanes  $\{x \in \mathbb{R}^n \mid a^\top x = b_1\}$  and  $\{x \in \mathbb{R}^n \mid a^\top x = b_2\}$ ? Hint. Let  $a^\top x_1 = b_1$ ,  $a^\top x_2 = b_2$ , and minimize  $\|x_1 - x_2\|^2$ .

Awr: Hyperplanes are given as =

$$\text{hyperplane 1} = \{x \in \mathbb{R}^n \mid a^\top x = b_1\}$$

$$\text{hyperplane 2} = \{x \in \mathbb{R}^n \mid a^\top x = b_2\}$$

let  $x_1$  be any point in the front hyperplane and consider that there is a line ' $l$ ' which is passing through  $x_1$  in the direction of the normal vector  $a$ .

Since point  $x_1$  is on hyperplane and line  $l$  is passing through it we can write the equation of line  $l$  as =

$$l = x_1 + at, t \in \mathbb{R}$$

the line is intersecting the hyperplane 2. Hence intersection of hyperplane 2 with line  $l$  is given as =

$$a^T(x_1 + a\tau) = b_2$$

$$\Rightarrow a^T x_1 + a^T a \tau = b_2$$

$$\Rightarrow a^T a \tau = b_2 - a^T x_1$$

$$\Rightarrow \tau = \frac{b_2 - a^T x_1}{a^T a}$$

But we know that  $\boxed{a^T x_1 = b_1}$  Hence all  
can write  $\tau =$

$$\boxed{\tau = \frac{b_2 - b_1}{a^T a}}$$

Therefore the intersection point on  
hyperplane 2 will be given as

$$\boxed{x_2 = x_1 + a \frac{(b_2 - b_1)}{a^T a}}$$

The distance between the points  $x_1, x_2$  is the distance between the hyperplanes. From above equation we can obtain =

$$x_2 - x_1 = \frac{a(b_2 - b_1)}{a^T a}$$

taking norm to find distance =

$$\|x_1 - x_2\| = \frac{\|a\|(b_2 - b_1)}{a^T a}$$

$$\boxed{\|x_1 - x_2\| = \frac{|b_2 - b_1|}{\|a\|}}$$

2. (20 points) Let  $x$  be a real-valued random variable with sample space  $\{a_1, \dots, a_k\}$  where  $a_1 \leq a_2 \leq \dots \leq a_k$ . This can be viewed as a categorical random variable with each category assigned a real value. Let  $\Pr[x = a_i] = p_i$ , then the vector  $\mathbf{p}$  satisfies  $\mathbf{p} \geq 0$  and  $\mathbf{1}^\top \mathbf{p} = 1$ , i.e., it lies in the probability simplex  $\Delta$ . For each of the following functions of  $\mathbf{p}$  on the probability simplex, determine if the function is convex, concave, or neither.

- (a)  $E[x]$
- (b)  $\Pr[x > \alpha]$
- (c)  $\Pr[\alpha < x < \beta]$
- (d)  $-\sum_{i=1}^k p_i \log p_i$ , the entropy of this distribution
- (e)  $\text{var}(x)$

Ans-2 =  $x$  is a real valued random variable with sample space  $\{a_1, a_2, \dots, a_k\}$ . where  $(a_1 \leq a_2 \leq \dots \leq a_k)$ .

Also given that  $\Pr[x = a_i] = p_i$   
then the vector  $\mathbf{p}$  satisfies  $\mathbf{p} \geq 0$  and  $\mathbf{1}^\top \mathbf{p} = 1$

(Part a) =

$$E(x) = a_1 p_1 + a_2 p_2 + \dots + a_k p_k = \mathbf{a}^\top \mathbf{p}$$

This function is a linear function.

Hence it is both concave & convex.

(Part b) =

$$\Pr[x > \alpha] = p_{\alpha+1} + p_{\alpha+2} + \dots + p_K$$

$$= L^T p$$

Here  $L$  starts from  $(\alpha+1)$  to  $K$  and all terms before  $(\alpha+1)$  are 0.

again this is also a linear function  
Hence it is both convex and concave.

(Part c) =

$$\Pr[\alpha < x < \beta] = p_{\alpha+1} + p_{\alpha+2} + \dots + p_\beta$$

$$= L^T p$$

Just like (Part b), Here also we can interpret

$$\Pr[\alpha < x < \beta] \text{ or } L^T p.$$

$$\text{where } L[0 : \alpha] = 0$$

$$L[\alpha+1 : \beta] = 1$$

$$L[\beta : K] = 0$$

Since this function is also linear function  
Hence it is also both concave & convex.

$$\text{(Part d)} = - \sum_{i=1}^k p_i \log p_i$$

we know that  $\log x$  is concave function.

But  $x \log x$  is a convex function.

Hence  $\sum_{i=1}^k p_i \log p_i$  is also a convex function.

But we know that  $\rightarrow$

if  $f(x) \rightarrow$  convex

Then  $-f(x) \rightarrow$  concave

Hence  $- \sum_{i=1}^k p_i \log p_i$  is a concave function

(Part e) =

we know that, (variance)  $\Rightarrow$

$$\text{var}(x) = E_p[x^2] - E_p[x]^2$$

also for  $\lambda \in [0,1]$  we can write that =

$$E_{\lambda p + (1-\lambda)q} [x^2] = \lambda E_p [x^2] + (1-\lambda) E_q [x^2]$$

$x^2$  is a convex function also the Jensen's Inequality gives us =

$$\begin{aligned} \lambda E_p [x]^2 + (1-\lambda) E_q [x]^2 &\geq (\lambda E_p [x] + (1-\lambda) E_q [x])^2 \\ &= E_{\lambda p + (1-\lambda)q} [x]^2 \end{aligned}$$

Thus we can infer =

$$\begin{aligned} \text{var}_{\lambda p + (1-\lambda)q}(x) &= E_{\lambda p + (1-\lambda)q} [x^2] - E_{\lambda p + (1-\lambda)q} [x]^2 \\ &\geq \lambda \text{var}_p(x) + (1-\lambda) \text{var}_q(x) \end{aligned}$$

This clearly indicates that variance  $\text{var}(x)$  is concave.

3. (10 points) Show that the following two convex problems are equivalent. Carefully explain how the solution of (b) is obtained from the solution of (a).

(a) The robust least squares problem

$$\underset{\boldsymbol{\theta}}{\text{minimize}} \quad \sum_{i=1}^n h(\boldsymbol{\phi}_i^\top \boldsymbol{\theta} - \psi_i),$$

where  $h : \mathbb{R} \rightarrow \mathbb{R}$  is the Huber function defined (with a constant  $M$ ) as

$$h(t) = \begin{cases} t^2 & |t| \leq M \\ M(2|t| - M) & |t| > M. \end{cases}$$

(b) The quadratic program

$$\begin{aligned} & \underset{\boldsymbol{\theta}, \mathbf{u}, \mathbf{v}}{\text{minimize}} \quad \sum_{i=1}^n (u_i^2 + 2Mv_i) \\ & \text{subject to} \quad -\mathbf{u} - \mathbf{v} \leq \boldsymbol{\Phi}\boldsymbol{\theta} - \boldsymbol{\psi} \leq \mathbf{u} + \mathbf{v} \\ & \quad 0 \leq \mathbf{u} \leq M\mathbf{1}, \quad \mathbf{v} \geq 0. \end{aligned}$$

Anw-3 = we have got,

robust least square problem  $\Rightarrow$

$$\underset{\boldsymbol{\theta}}{\text{minimize}} \sum_{i=1}^n h(\boldsymbol{\phi}_i^\top \boldsymbol{\theta} - \psi_i)$$

where  $h(t)$  is huber function

and the quadratic program =

$$\underset{\boldsymbol{\theta}, \mathbf{u}, \mathbf{v}}{\text{minimize}} \sum_{i=1}^n (u_i^2 + 2Mv_i)$$

which is subject to  $-\mathbf{u} - \mathbf{v} \leq \boldsymbol{\Phi}\boldsymbol{\theta} - \boldsymbol{\psi} \leq \mathbf{u} + \mathbf{v}$   
and  $0 \leq \mathbf{u} \leq M\mathbf{1}, \mathbf{v} \geq 0.$

we are given that  $\phi^T \theta - \psi \leq u + v$

we can rewrite it as =

$$|\phi_i^T \theta - \psi_i| \leq u_i + v_i$$

where  $u_i, v_i > 0$ , and  $u_i$  and  $v_i$  has to be minimized.

To minimize  $v_i$  we can write  $v_i$  as =

$$\boxed{v_i = |\phi_i^T \theta - \psi_i| - u_i} \rightarrow \textcircled{1}$$

lets substitute this value to the Quadratic Program =

$$\underset{\theta, u}{\text{minimize}} \quad \sum_{i=1}^n (u_i^2 + 2M(|\phi_i^T \theta - \psi_i| - u_i))$$

$$= \sum_{i=1}^n (u_i^2 + 2M|L_i| - 2Mu_i) \rightarrow \textcircled{2}$$

if we assume that  $L_i = |\phi_i^T \theta - \psi_i|$

$$\text{Then we have } u_i = \begin{cases} |L_i| & \text{when } |L_i| \leq M \\ M & |L_i| > M \end{cases}$$

$$\text{using equation } u_i = \begin{cases} |L_i| & \text{when } |L_i| \leq M \\ M & |L_i| > M \end{cases}$$

we can modify equation ②.

$$\textcircled{1} = \text{when } |L_i| \leq M \text{ then } u_i = |L_i|$$

so equation ② becomes =

$$((|L_i|)^2 + 2M|L_i| - 2M|L_i|) = L$$

(when  $|L| \leq M$ )

$$\textcircled{2} =$$

also when  $|L_i| > M$  then  $u_i = M$

so equation ② becomes =

$$(M^2 - 2M|L_i| - 2M^2) = M(2|L| - M)$$

(when  $|L| > M$ )

by writing above two cases together  
we can write =

$$\text{minimize}_{\theta} f(\theta) = \begin{cases} L & |L| \leq M \\ M(2|L|-M) & |L| > M \end{cases}$$

(Huber loss)  
function

which is name as saying =

$$\text{minimize}_{\theta} \sum_{i=1}^n h(L_i) \quad [L_i = |\phi_i^\top \theta - \psi_i|]$$

$$\Rightarrow \text{minimize}_{\theta} \sum_{i=1}^n h(\phi_i^\top \theta - \psi_i)$$

Hence both of these equations in part (a) and (b) are equivalent. and so (b) can be obtained by solution of (a).

4. (30 points) We test the performance of three regression methods on the wine data set <http://archive.ics.uci.edu/ml/datasets/Wine+Quality>. We will only consider the red wine data set, with 1599 samples. We use the first 1400 samples for training, and the last 199 samples for testing. The goal is to build a linear model of the first 11 features (together with a constant term) to predict the quality of the wine. All models are trained by solving the following optimization problem

$$\underset{\mathbf{w}, \beta}{\text{minimize}} \quad \sum_{i=1}^n \ell(\mathbf{x}_i^\top \mathbf{w} + \beta - y_i),$$

where the loss functions are

- least squares loss  $\ell(t) = t^2$
- Huber loss defined in the previous problem, with  $M = 1$
- hinge (deadzone-linear) loss

$$\ell(t) = \begin{cases} 0 & |t| \leq 0.5 \\ |t| - 0.5 & |t| > 0.5 \end{cases}$$

The least squares loss can be directly solved by the command `Phi\y` for some properly defined `Phi`. For the latter two, you will use the `cvx` package found on Prof. Boyd's website <https://web.stanford.edu/~boyd/software.html>. Report their prediction performance on the test set using a different metric, mean absolute error (MAE), defined as  $(1/n) \sum_{i=1}^n |y_i - \hat{y}_i|$ .

Ans-4 = I have submitted the working code for this question along with this assignment. Please check the attached files.

(a) = least squares loss  $\ell(t) = t^2$

MAE using phi matrix = 0.674940

Prepared phi matrix using various quadratic basis functions for 11 features. Please check code & comments for more info.

(b) = Huber loss function with  $\gamma = 1$   
MAE using Huber loss = 0.53271  
(using CVX package)

(c) = Hinge loss function  $\Rightarrow$

MAE using Hinge loss = 0.54810  
(using CVX package)

among these three functions for loss my model is giving least MAE using Huber loss.



Note = please check code for all phi matrix, CVX loss function

Implementation.

② = I have implemented least squares using phi matrix & CVX package both but here I am giving phi-matrix results.

## Model training & testing Glossary =

- ① = available data was divided into training & test data sets. each set had feature vectors data & outcome data separately ( $\gamma$ -train,  $\gamma$ -train,  $X$ -test,  $\gamma$ -test)
- ② = least squares Problem was solved by forming  $\phi$  matrix using quadratic basis functions for features. quadratic function was used b/c it was accurate at the same time it was generalizing. for  $\phi^+$  = Pseudo inverse was calculated
- $$(\phi^\top \phi)^{-1} \phi^\top$$
- Then  $\theta = \phi^+ (\gamma\text{-train})$

after getting  $\theta$  parameters. predictions were made on test data, these prediction values were fed to activation function (Sigmoid) to obtain classes.

③ = Huber loss & Hinge loss both were optimized using CVX packages. and parameters in  $(\mathbf{z}\mathbf{w} + \beta)$  were obtained.  $(\mathbf{w}, \beta)$ . using these parameter prediction on test data was made. these prediction values were fed to activation function (Sigmoid) to obtain classes.

④ = after making predictions MAE was calculated.

5. (30 points) We test the performance of three classification methods on the ionosphere data set <https://archive.ics.uci.edu/ml/datasets/ionosphere>. There are 351 samples. We use the first 300 samples for training, and the last 51 samples for testing. The goal is to build a linear model of the 34 features (together with a constant term) to predict the binary ( $\pm 1$ ) outcome. All models are trained by solving the following optimization problem

$$\underset{\mathbf{w}, \beta}{\text{minimize}} \quad \sum_{i=1}^n \ell(\mathbf{x}_i^\top \mathbf{w} + \beta, y_i),$$

where the loss functions are

- least squares loss  $\ell(t, y) = (t - y)^2$
- logistic loss  $\ell(t, y) = \log(1 + \exp(-yt))$
- hinge loss  $\ell(t, y) = \max(0, 1 - yt)$

Again, you will use the backslash command to solve for the first model, and `cvx` to solve for the latter two. Report their prediction accuracy on the test set.

Ans-4 = I have submitted the working code for this question along with this assignment. Please check the attached files.

(1) = Least Squares loss  $\Rightarrow$  (using  $\Phi$  matrix)  
Classifier model Accuracy = 90.19 %

(2) = Logistic loss  $\Rightarrow$  (using CVX package)  
Classifier model Accuracy = 74.50 %

(3) = Hinge loss = (using CVX package)  
Classifier model accuracy = 74.50 %

30

least square model accuracy  $\rightarrow$  logistic model accuracy = Hinge loss model accuracy

Note = please check code for all phi matrix, CVX loss function implementation.

- ② = I have implemented least squares loss using phi matrix & CVX package both but here I have shared the result I obtained using only phi matrix method.
- ③ = ideally (Hinge loss, logistic loss) should not give equal accuracy but that's what I obtained even after tuning the model.

## Model training & testing Glossary =

- ① = available data was divided into training & test datasets. each set had feature vectors data & outcome data separately ( $\gamma$ -train,  $\gamma$ -test,  $X$ -train,  $X$ -test)
- ② = least squares Problem was solved by forming  $\phi$  matrix using quadratic basis functions for features. quadratic function was used b/c it was accurate at the same time it was generalizing.  
for  $\phi^+ = \text{Pseudo inverse was calculated}$   
$$(\phi^\top \phi)^{-1} \phi^\top$$

$$\text{Then } \theta = \phi^+ (\gamma\text{-train})$$

after getting  $\theta$  parameters. predictions were made on test data, these prediction values were fed to activation function (Sigmoid) to obtain classes.

③= Logistic loss & Hinge loss both were optimized using CVX packages. and parameters in  $(\mathbf{x}^T \mathbf{w} + \beta)$  were obtained.  $(\mathbf{w}, \beta)$ . Using these parameter prediction on test data was made. These prediction values were fed to activation function (Sigmoid) to obtain classes.

④= Accuracy was calculated by comparing available actual outcomes with Predictions.